

## IP VIRTUALIZATION

### REFERENCE TO RELATED APPLICATIONS

The present application claims priority to and incorporates the following applications by reference: DYNAMIC SYMBOLIC LINK RESOLUTION, Prov. No. 60/157,728, filed on October 5, 1999; SNAPSHOT VIRTUAL TEMPLATING, Prov. No. 60/157,728, filed on October 5, 1999; SNAPSHOT RESTORE OF APPLICATION CHAINS AND APPLICATIONS, Prov. No. 60/157,833, filed on October 5, 1999; VIRTUAL RESOURCE-ID MAPPING, Prov. No. 60/157,727, filed on October 5, 1999; and VIRTUAL PORT MULTIPLEXING, Prov. No. 60/157,834, filed on October 5, 1999.

### FIELD

The present invention relates broadly to client server computer systems and networks. More specifically, the present invention relates to virtualization of a network identity to allow application migration across computers in a computer network.

### BACKGROUND

Global computer networks such as the Internet have allowed electronic commerce ("e-commerce") to flourish to a point where a large number of customers purchase goods and services over websites operated by online merchants. Because the Internet provides an effective medium to reach this large customer base, online merchants who are new to the e-commerce marketplace are often flooded with high customer traffic from the moment their websites are rolled out. In order to effectively serve customers, online merchants are charge with the same responsibility as conventional merchants: they must provide quality service to their customers in a timely manner. Often, insufficient computing resources are the cause of a bottleneck that results in customer frustration and loss of sales. This phenomena has resulted in the need for a new utility: leasable online computing infrastructure. Previous attempts at providing computing resources have entailed leasing large blocks of storage and processing power. However, for a new online merchant having no baseline from which to judge customer traffic upon rollout, this

approach is inefficient. Either too much computing resources are leased, depriving a start up merchant of financial resources that are needed elsewhere in the operation, or not enough resources are leased, and a bottleneck occurs. Thus, there remains a heartfelt need for computer infrastructure that can be provided in an on demand basis.

5 One impediment to providing an on-demand computer infrastructure involves application migration. A migrating application instance, along with any applications the migrating instance is communicating with, must maintain a consistent view of the migrating application's network identity for interoperability. For example, consider a client application that is communicating with a server application having an IP address of a.b.c.d. The client application initiates two  
10 sequential data connections to a server. Between the first and the second connection, however, the server gets moved to a network node having the IP address a.b.c.z. If the IP address of the server was a real IP address (associated with a physical interface) of the first host, the second connection would not address the correct machine.

## 15 SUMMARY

The present invention solves the problems described above by providing virtualization of network parameters to create a virtual network identity for an application. The virtualization of the network parameters, IP address and hostname, allows for the migration of an application instance. By virtualizing these network parameters, a virtual network identity (VNI) is created for the application instance. This virtual network identity remains with the application regardless of which node the application is running on. This allows the application to have a static network identity. Associating the IP address and hostname with a running instance of an application distinguishes the present invention from the prior art where an IP address is associated with a network interface, either physical or loopback, of an individual computer. With the virtual IP  
20 address method and system of the present invention, the IP address and hostname are associated with a particular application instance. Only the processes within the application instance can use the virtual IP address assigned to it. Furthermore, the application instance can only use the virtual address, binding the application to the VNI by preventing the application from using other (real or virtual) addresses that may be resident in the host.

## BRIEF DESCRIPTION OF THE DRAWINGS

FIG. 1 is a high level block diagram illustrating the various components of a computer network used in connection with the present invention;

5 FIG. 2 is a high level block diagram illustrating the various components of a computer used in connection with the present invention;

FIG. 3 illustrates how application state is tracked using library and operating system kernel interposition;

FIG. 4 illustrates the capture of an application's run-time state;

10 FIG. 5 is a flow chart illustrating the logical sequence of steps executed to create a snapshot image of an application instance;

FIG. 6 is a flow chart illustrating the logical sequence of steps executed to restore an application instance from a snapshot image.;

FIG. 7 is an illustration of the format of a snapshot virtual template;

15 FIG. 8 is a flowchart illustrating the logical sequence of steps executed to create a snapshot virtual template;

FIG. 9 is a flowchart illustrating the logical sequence of steps executed to clone a snapshot virtual template;

FIG. 10 illustrates the registration of an application using virtual resource identifiers;

FIG. 11 illustrates the allocation of a virtual resource;

20 FIG. 12 illustrates the translation of a virtual resource to a system resource;

FIG. 13 illustrates the translation of a system resource to a virtual resource;

FIG. 14 is a flowchart illustrating the logical sequence of steps executed to create a virtual translation table;

25 FIG. 15 is a flowchart illustrating the logical sequence of steps executed to translate a virtual resource;

FIG. 16 illustrates the resolution of the local hostname for an application running with the VNI;

FIG. 17 illustrates how a server application running within a VNI is connected to by a client application; and

30 FIG. 18 illustrates how a client application running within a VNI connects with a server application.

## DETAILED DESCRIPTION

### A. Snapshot Restore

FIG. 1 illustrates in high level block diagram form the overall structure of the present invention as used in connection with a global computer network 100 such as the Internet.

5 Remote users 102-1 and 102-2 can connect through the computer network 100 to a private network of computers 106 protected by firewall 104. Computer network 106 is a network comprising computers 150-1, 150-2, through 150-n, where n is the total number of computers in network 106. Computers 150 are used to run various applications, as well as host web sites for access by remote users 102. The present invention is implemented on computer network 106 in 10 the form of virtual environments 110-1 and 110-2. While only two virtual environments are illustrated, it is to be understood that any number of virtual environments may be utilized in connection with the present invention.

0 5 10 15 20 25 30 35 40 45 50 55 60 65 70 75 80 85 90 95 100 105 110 115 120 125 130 135 140 145 150 155 160 165 170 175 180 185 190 195 200 205 210 215 220 225 230 235 240 245 250 255 260 265 270 275 280 285 290 295 300 305 310 315 320 325 330 335 340 345 350 355 360 365 370 375 380 385 390 395 400 405 410 415 420 425 430 435 440 445 450 455 460 465 470 475 480 485 490 495 500 505 510 515 520 525 530 535 540 545 550 555 560 565 570 575 580 585 590 595 600 605 610 615 620 625 630 635 640 645 650 655 660 665 670 675 680 685 690 695 700 705 710 715 720 725 730 735 740 745 750 755 760 765 770 775 780 785 790 795 800 805 810 815 820 825 830 835 840 845 850 855 860 865 870 875 880 885 890 895 900 905 910 915 920 925 930 935 940 945 950 955 960 965 970 975 980 985 990 995 1000 1005 1010 1015 1020 1025 1030 1035 1040 1045 1050 1055 1060 1065 1070 1075 1080 1085 1090 1095 1100 1105 1110 1115 1120 1125 1130 1135 1140 1145 1150 1155 1160 1165 1170 1175 1180 1185 1190 1195 1200 1205 1210 1215 1220 1225 1230 1235 1240 1245 1250 1255 1260 1265 1270 1275 1280 1285 1290 1295 1300 1305 1310 1315 1320 1325 1330 1335 1340 1345 1350 1355 1360 1365 1370 1375 1380 1385 1390 1395 1400 1405 1410 1415 1420 1425 1430 1435 1440 1445 1450 1455 1460 1465 1470 1475 1480 1485 1490 1495 1500 1505 1510 1515 1520 1525 1530 1535 1540 1545 1550 1555 1560 1565 1570 1575 1580 1585 1590 1595 1600 1605 1610 1615 1620 1625 1630 1635 1640 1645 1650 1655 1660 1665 1670 1675 1680 1685 1690 1695 1700 1705 1710 1715 1720 1725 1730 1735 1740 1745 1750 1755 1760 1765 1770 1775 1780 1785 1790 1795 1800 1805 1810 1815 1820 1825 1830 1835 1840 1845 1850 1855 1860 1865 1870 1875 1880 1885 1890 1895 1900 1905 1910 1915 1920 1925 1930 1935 1940 1945 1950 1955 1960 1965 1970 1975 1980 1985 1990 1995 2000 2005 2010 2015 2020 2025 2030 2035 2040 2045 2050 2055 2060 2065 2070 2075 2080 2085 2090 2095 2100 2105 2110 2115 2120 2125 2130 2135 2140 2145 2150 2155 2160 2165 2170 2175 2180 2185 2190 2195 2200 2205 2210 2215 2220 2225 2230 2235 2240 2245 2250 2255 2260 2265 2270 2275 2280 2285 2290 2295 2300 2305 2310 2315 2320 2325 2330 2335 2340 2345 2350 2355 2360 2365 2370 2375 2380 2385 2390 2395 2400 2405 2410 2415 2420 2425 2430 2435 2440 2445 2450 2455 2460 2465 2470 2475 2480 2485 2490 2495 2500 2505 2510 2515 2520 2525 2530 2535 2540 2545 2550 2555 2560 2565 2570 2575 2580 2585 2590 2595 2600 2605 2610 2615 2620 2625 2630 2635 2640 2645 2650 2655 2660 2665 2670 2675 2680 2685 2690 2695 2700 2705 2710 2715 2720 2725 2730 2735 2740 2745 2750 2755 2760 2765 2770 2775 2780 2785 2790 2795 2800 2805 2810 2815 2820 2825 2830 2835 2840 2845 2850 2855 2860 2865 2870 2875 2880 2885 2890 2895 2900 2905 2910 2915 2920 2925 2930 2935 2940 2945 2950 2955 2960 2965 2970 2975 2980 2985 2990 2995 3000 3005 3010 3015 3020 3025 3030 3035 3040 3045 3050 3055 3060 3065 3070 3075 3080 3085 3090 3095 3100 3105 3110 3115 3120 3125 3130 3135 3140 3145 3150 3155 3160 3165 3170 3175 3180 3185 3190 3195 3200 3205 3210 3215 3220 3225 3230 3235 3240 3245 3250 3255 3260 3265 3270 3275 3280 3285 3290 3295 3300 3305 3310 3315 3320 3325 3330 3335 3340 3345 3350 3355 3360 3365 3370 3375 3380 3385 3390 3395 3400 3405 3410 3415 3420 3425 3430 3435 3440 3445 3450 3455 3460 3465 3470 3475 3480 3485 3490 3495 3500 3505 3510 3515 3520 3525 3530 3535 3540 3545 3550 3555 3560 3565 3570 3575 3580 3585 3590 3595 3600 3605 3610 3615 3620 3625 3630 3635 3640 3645 3650 3655 3660 3665 3670 3675 3680 3685 3690 3695 3700 3705 3710 3715 3720 3725 3730 3735 3740 3745 3750 3755 3760 3765 3770 3775 3780 3785 3790 3795 3800 3805 3810 3815 3820 3825 3830 3835 3840 3845 3850 3855 3860 3865 3870 3875 3880 3885 3890 3895 3900 3905 3910 3915 3920 3925 3930 3935 3940 3945 3950 3955 3960 3965 3970 3975 3980 3985 3990 3995 4000 4005 4010 4015 4020 4025 4030 4035 4040 4045 4050 4055 4060 4065 4070 4075 4080 4085 4090 4095 4100 4105 4110 4115 4120 4125 4130 4135 4140 4145 4150 4155 4160 4165 4170 4175 4180 4185 4190 4195 4200 4205 4210 4215 4220 4225 4230 4235 4240 4245 4250 4255 4260 4265 4270 4275 4280 4285 4290 4295 4300 4305 4310 4315 4320 4325 4330 4335 4340 4345 4350 4355 4360 4365 4370 4375 4380 4385 4390 4395 4400 4405 4410 4415 4420 4425 4430 4435 4440 4445 4450 4455 4460 4465 4470 4475 4480 4485 4490 4495 4500 4505 4510 4515 4520 4525 4530 4535 4540 4545 4550 4555 4560 4565 4570 4575 4580 4585 4590 4595 4600 4605 4610 4615 4620 4625 4630 4635 4640 4645 4650 4655 4660 4665 4670 4675 4680 4685 4690 4695 4700 4705 4710 4715 4720 4725 4730 4735 4740 4745 4750 4755 4760 4765 4770 4775 4780 4785 4790 4795 4800 4805 4810 4815 4820 4825 4830 4835 4840 4845 4850 4855 4860 4865 4870 4875 4880 4885 4890 4895 4900 4905 4910 4915 4920 4925 4930 4935 4940 4945 4950 4955 4960 4965 4970 4975 4980 4985 4990 4995 5000 5005 5010 5015 5020 5025 5030 5035 5040 5045 5050 5055 5060 5065 5070 5075 5080 5085 5090 5095 5100 5105 5110 5115 5120 5125 5130 5135 5140 5145 5150 5155 5160 5165 5170 5175 5180 5185 5190 5195 5200 5205 5210 5215 5220 5225 5230 5235 5240 5245 5250 5255 5260 5265 5270 5275 5280 5285 5290 5295 5300 5305 5310 5315 5320 5325 5330 5335 5340 5345 5350 5355 5360 5365 5370 5375 5380 5385 5390 5395 5400 5405 5410 5415 5420 5425 5430 5435 5440 5445 5450 5455 5460 5465 5470 5475 5480 5485 5490 5495 5500 5505 5510 5515 5520 5525 5530 5535 5540 5545 5550 5555 5560 5565 5570 5575 5580 5585 5590 5595 5600 5605 5610 5615 5620 5625 5630 5635 5640 5645 5650 5655 5660 5665 5670 5675 5680 5685 5690 5695 5700 5705 5710 5715 5720 5725 5730 5735 5740 5745 5750 5755 5760 5765 5770 5775 5780 5785 5790 5795 5800 5805 5810 5815 5820 5825 5830 5835 5840 5845 5850 5855 5860 5865 5870 5875 5880 5885 5890 5895 5900 5905 5910 5915 5920 5925 5930 5935 5940 5945 5950 5955 5960 5965 5970 5975 5980 5985 5990 5995 6000 6005 6010 6015 6020 6025 6030 6035 6040 6045 6050 6055 6060 6065 6070 6075 6080 6085 6090 6095 6100 6105 6110 6115 6120 6125 6130 6135 6140 6145 6150 6155 6160 6165 6170 6175 6180 6185 6190 6195 6200 6205 6210 6215 6220 6225 6230 6235 6240 6245 6250 6255 6260 6265 6270 6275 6280 6285 6290 6295 6300 6305 6310 6315 6320 6325 6330 6335 6340 6345 6350 6355 6360 6365 6370 6375 6380 6385 6390 6395 6400 6405 6410 6415 6420 6425 6430 6435 6440 6445 6450 6455 6460 6465 6470 6475 6480 6485 6490 6495 6500 6505 6510 6515 6520 6525 6530 6535 6540 6545 6550 6555 6560 6565 6570 6575 6580 6585 6590 6595 6600 6605 6610 6615 6620 6625 6630 6635 6640 6645 6650 6655 6660 6665 6670 6675 6680 6685 6690 6695 6700 6705 6710 6715 6720 6725 6730 6735 6740 6745 6750 6755 6760 6765 6770 6775 6780 6785 6790 6795 6800 6805 6810 6815 6820 6825 6830 6835 6840 6845 6850 6855 6860 6865 6870 6875 6880 6885 6890 6895 6900 6905 6910 6915 6920 6925 6930 6935 6940 6945 6950 6955 6960 6965 6970 6975 6980 6985 6990 6995 7000 7005 7010 7015 7020 7025 7030 7035 7040 7045 7050 7055 7060 7065 7070 7075 7080 7085 7090 7095 7100 7105 7110 7115 7120 7125 7130 7135 7140 7145 7150 7155 7160 7165 7170 7175 7180 7185 7190 7195 7200 7205 7210 7215 7220 7225 7230 7235 7240 7245 7250 7255 7260 7265 7270 7275 7280 7285 7290 7295 7300 7305 7310 7315 7320 7325 7330 7335 7340 7345 7350 7355 7360 7365 7370 7375 7380 7385 7390 7395 7400 7405 7410 7415 7420 7425 7430 7435 7440 7445 7450 7455 7460 7465 7470 7475 7480 7485 7490 7495 7500 7505 7510 7515 7520 7525 7530 7535 7540 7545 7550 7555 7560 7565 7570 7575 7580 7585 7590 7595 7600 7605 7610 7615 7620 7625 7630 7635 7640 7645 7650 7655 7660 7665 7670 7675 7680 7685 7690 7695 7700 7705 7710 7715 7720 7725 7730 7735 7740 7745 7750 7755 7760 7765 7770 7775 7780 7785 7790 7795 7800 7805 7810 7815 7820 7825 7830 7835 7840 7845 7850 7855 7860 7865 7870 7875 7880 7885 7890 7895 7900 7905 7910 7915 7920 7925 7930 7935 7940 7945 7950 7955 7960 7965 7970 7975 7980 7985 7990 7995 8000 8005 8010 8015 8020 8025 8030 8035 8040 8045 8050 8055 8060 8065 8070 8075 8080 8085 8090 8095 8100 8105 8110 8115 8120 8125 8130 8135 8140 8145 8150 8155 8160 8165 8170 8175 8180 8185 8190 8195 8200 8205 8210 8215 8220 8225 8230 8235 8240 8245 8250 8255 8260 8265 8270 8275 8280 8285 8290 8295 8300 8305 8310 8315 8320 8325 8330 8335 8340 8345 8350 8355 8360 8365 8370 8375 8380 8385 8390 8395 8400 8405 8410 8415 8420 8425 8430 8435 8440 8445 8450 8455 8460 8465 8470 8475 8480 8485 8490 8495 8500 8505 8510 8515 8520 8525 8530 8535 8540 8545 8550 8555 8560 8565 8570 8575 8580 8585 8590 8595 8600 8605 8610 8615 8620 8625 8630 8635 8640 8645 8650 8655 8660 8665 8670 8675 8680 8685 8690 8695 8700 8705 8710 8715 8720 8725 8730 8735 8740 8745 8750 8755 8760 8765 8770 8775 8780 8785 8790 8795 8800 8805 8810 8815 8820 8825 8830 8835 8840 8845 8850 8855 8860 8865 8870 8875 8880 8885 8890 8895 8900 8905 8910 8915 8920 8925 8930 8935 8940 8945 8950 8955 8960 8965 8970 8975 8980 8985 8990 8995 9000 9005 9010 9015 9020 9025 9030 9035 9040 9045 9050 9055 9060 9065 9070 9075 9080 9085 9090 9095 9100 9105 9110 9115 9120 9125 9130 9135 9140 9145 9150 9155 9160 9165 9170 9175 9180 9185 9190 9195 9200 9205 9210 9215 9220 9225 9230 9235 9240 9245 9250 9255 9260 9265 9270 9275 9280 9285 9290 9295 9300 9305 9310 9315 9320 9325 9330 9335 9340 9345 9350 9355 9360 9365 9370 9375 9380 9385 9390 9395 9400 9405 9410 9415 9420 9425 9430 9435 9440 9445 9450 9455 9460 9465 9470 9475 9480 9485 9490 9495 9500 9505 9510 9515 9520 9525 9530 9535 9540 9545 9550 9555 9560 9565 9570 9575 9580 9585 9590 9595 9600 9605 9610 9615 9620 9625 9630 9635 9640 9645 9650 9655 9660 9665 9670 9675 9680 9685 9690 9695 9700 9705 9710 9715 9720 9725 9730 9735 9740 9745 9750 9755 9760 9765 9770 9775 9780 9785 9790 9795 9800 9805 9810 9815 9820 9825 9830 9835 9840 9845 9850 9855 9860 9865 9870 9875 9880 9885 9890 9895 9900 9905 9910 9915 9920 9925 9930 9935 9940 9945 9950 9955 9960 9965 9970 9975 9980 9985 9990 9995 9999

200 is transparent to running (and snapshotted) applications. From an application's perspective, the application is always running. An application snapshot may consist of multiple processes and multiple threads and includes shared resources in use by a process, such as shared memory or semaphores. A process may be snapshotted & restored more than once. The computer on which a process is restored on must be identically configured and have an identical environment (hardware, software, and files) that matches the environment of the computer where the process was snapshotted. All processes that are snapshotted together in the form of an application chain share the same application ID ("AID"). As used herein, an application chain is the logical grouping of a set of applications and processes that communicate with each other and share resources to provide a common function.

~~The virtual environment 204 is a layer that surrounds application(s) 208 and resides between the application and the operating system 206. Resource handles are abstracted to present a consistent view to the application although the actual system resource handles may change as an application is snapshot/restored more than once. The virtual environment also allows multiple applications to compete for the same resources where exclusion would normally prohibit such behavior to allow multiple snapshots to coexist without reconfiguration. Preload library 214 is an application library that interposes upon an application for the express purpose of intercepting and handling library called and system calls. Once the library has been preloaded it is attached to the process' address space. Preload library 214 interposes between application 208 and operating system 206. It is distinguished from kernel interposition in that it operates in "user mode" (i.e., non-kernel and non-privileged mode). Application 208 can make application programming interface (API) calls that modify the state of the application. These calls are made from the application 208 to the operating system API interfaces 210 via the application snapshot restore framework 200 or the preload library 214. The preload library can save the state of various resources by intercepting API interface calls and then saves the state at a pre-arranged memory location. When the process' memory is saved as part of the snapshot/restore mechanism, this state is saved since it resides in memory. The state as it is modified is saved to non-volatile storage (i.e. a file on disk). The preload library notify the snapshot/restore framework through one of its private interface.~~

30 FIG. 4 illustrates the capture of an application's run time state. The OS API interfaces

210 are standard programming interfaces defined by international standards organizations such as XOPEN. The open() system call which allows an application to open a file for reading is an example of an API interface. The process management system 216 is a component of the operating system 206 that allows one process to examine or alter the state of another process.

5 The interfaces that are provided by this component are usually not standardized interfaces (not part of a recognized standard API) and are OS-implementation dependent. However, such interfaces usually allow access to more state than standardized API interfaces. The run-time information captured from a process is used by the snapshot driver 218.

An application needs to be snapshotted if it is idle and is not currently in use or there are  
10 higher priority requests that require the application be scheduled out and preempted in favor of another application. A snapshot request is initiated by an application scheduler that determines when an application needs to be restored on-demand and when the application is no longer needed (can be snapshotted to free up resources). The application scheduler does this based on web traffic, server load, request response time, and a number of other factors. An application needs to be restored if there is an incoming request (i.e. a web browser request) and the application required to handle that request (ie a particular web site) is not currently running. Alternatively, an application needs to be restored if there is an incoming request (i.e. a web browser request) and the application required to handle that request (ie a particular web site) is currently overloaded, so another instance of that application is restored to handle that request.

200 FIG. 5 illustrates the logical sequence of steps executed by Snapshot driver 218 to make a snapshot image of a process. Beginning at step 250, an snapshot image of a runnable application is requested. The AID is looked up (decision step 252) in a table in memory 154 containing a list of every AID present on computer 150. If the AID is not found control returns at step 254. However, if the AID is found, control continues to decision step 256 where the snapshot/restore  
25 framework 200 searches for a process belonging to the application having the matched AID. If a process is found, control continues to step 258, where the process is suspended. For a process to be snapshotted, it must be completely suspended with no activity present and no ambiguous state (i.e., in a transitory state). Since a process may be represented by asynchronous threads of activity in the operating system that are not formally part of the process state, any activity in the  
30 operating system that is executing on behalf of the process must be stopped (i.e. disk I/O

activity). In other words, there may be moments where temporarily a process cannot be snapshotted. This is a finite and short period of time, but it can occur. If the state is consistent and the threads are quiesced (decision step 260), control loops to step 256 and the remaining processes belonging to the application are located and suspended. However, if a process is 5 located that does not have a consistent state or a thread is not quiesced, suspended processes are resumed and the snapshot cannot be completed.

Once all related processes are suspended, for each state of each suspended process, the state is checked to see if it is virtualized (step 264). A virtualized state is any process state that reflects a virtualized resource. If the state is virtualized, it is retrieved at step 266 ; otherwise the 10 non-virtualized state is retrieved at step 268 State retrieval is performed as described above by the snapshot driver 218 querying the application snapshot/restore framework 200, operating system API interfaces 210, and process management subsystem 216. If the state has changed since the last snapshot (step 270), the new state is recorded. Control then loops to step 264 and executes through the above sequence of steps until all states of all processes are checked. Once 15 completed, control proceeds to step 278, the registered global state, such as semaphores, is removed. Registered global state is state that is not specifically associated with any one process (ie private state). Global state is usually exported (accessible) to all processes and its state can be modified (shared) by all processes. Control proceeds to step 280, where the process is terminated. If there are remaining processes (step 282), these are also terminated. This sequence 20 of steps is concluded to create a snapshot image which is stored as a file and made available for transmission to another computer within public computer network 100 or private computer network 106.

FIG. 6 illustrates the sequence of steps executed by the restore driver 220 to restore a snapshot image. The snapshot image is accessed via a shared storage mechanism and a restore 25 call is made at step 300. The restore driver 220 looks up the AID for the snapshot image and (decision step 302) if not found control returns and the snapshot image cannot be restored. However, if the AID is found, control continues to decision step 304 where, if the snapshotted image matching the AID is located, the global/shared state for each process associated with the snapshot are found. Control then continues to step 308, where remaining global or shared state 30 for the processes are recreated. Since global and shared state is not associated with a single

process and may be referenced by multiple processes, it is created first. Recreating the state entails creating a global resource that is functionally identical to the resource at the time of the snapshot. For example if during a snapshot, a semaphore with id 5193 is found with a value of 7, then to recreate the state at restore time a new semaphore must be created having the exact same 5 ID as before (ie 5193) and it also must have the same state (ie value 7) as before. Then, for each image, a process is created that inherits the global/shared state restored in step 308, and each created process is isolated to prevent inter-process state changes. When a process is being restored, process state is being registered with the kernel, inter-process mechanisms are being restored and reconnected and I/O buffers in the kernel may be being restored. Some of these 10 actions in one process may have the unintended side effect of disturbing another process that is also being restored. For example if an I/O buffer that is in the operating system as a result of a process, performing a write to a socket connection, then process, could unintentionally be delivered an asynchronous signal that notifies it of I/O being present (for reading) prior to the process being fully restored. At step 314, for each type of state within the processes, the process- 15 private resources are recreated to their state at the time the snapshot image was taken. If the state is virtualized (decision step 316), the system state is bound to a virtual definition. As part of the restore an extra step must be done to create a virtual mapping. This is done by taking the system resource that was created in step 314 and binding it to the virtual definition that was saved during the snapshot in step 266. This allows the application to see a consistent view of resources, since it cannot be guaranteed that at restore time the exact same system resource will 20 be available. If the state is shared with another process, such as via a pipe (decision state 320), the shared state is reconnected with the other process at step 322. If there are more states (decision step 324) steps 314 through 322 are repeated. Once steps 314 through 322 have been 25 executed for all states, control continues to step 326, where the process is placed in synchronized wait. If there are remaining images in the snapshot image (decision step 328), steps 310 through 326 are repeated. Once all images have been processed, control continues to step 330, where traces and states induced during restore of the process are removed, and a synchronized resume of all processes occurs at step 332.

Once steps 300 through 332 have executed without error on whatever computer the 30 restore driver 220 was executed, the restored application can continue to run without

interruption. Thus, the present invention avoids the overhead and delay of shutting down an application, storing data to a separate file, moving both the application and data file elsewhere, and restarting the program.

5        B.     Snapshot Virtual Templating

In another aspect, the present invention provides a system, method, and computer program product for creating snapshot virtual application templates for the purpose of propagating a single application snapshot into multiple distinct instances. Snapshot virtual templates allow multiple application instances to use the same fixed resource ID (“RID”) by 10 making the resource ID virtual, privatizing the virtual RID, and dynamically mapping it to a unique system resource ID. A RID is the identifier assigned to represent a specific system resource and acts as a handle when referencing that system resource. Anonymous resources are resources that are read-only or functionally isolated from other applications. Anonymous resources are also shareable resources. An anonymous resource is a non-fixed resource allocated 15 by the operating system and identified by a per-process handle. These are functionally-isolated since the operating system allocates it anonymously and one is as good as another. Examples of this are non-fixed TCP ports or file descriptors. A resource is said to be network-wide unique if there can only be one instance of that resource with its corresponding identifier on computer network or subnetwork. An example of this is a network IP address (i.e. 10.1.1.1). Snapshot virtual templates allow snapshots to be described in a manner that separates shareable data from 20 non-shareable data. Data is loosely defined to mean any system resource (memory, files, sockets, handles, etc.). When a snapshot is cloned from a virtual template, the common or shared data is used exactly as is, whereas the non-shareable data is either copied-on-write, multiplexed, virtualized, or customized-on-duplication. The present invention greatly reduces the required 25 administrative setup per application instance. Snapshot virtual templating works by noting access to modified resources, fixed system IDs/keys and unique process-related identifiers and automatically inserting a level of abstraction between these resources and the application. The resources contained in a snapshot virtual template can be dynamically redirected at restore time. Access to memory and storage is managed in a copy-on-write fashion. System resource handles 30 are managed in a virtualize-on-allocate fashion or by a multiplex-on-access mechanism. Process-

unique resources are managed in a redirect-on-duplicate fashion. Rules may be defined through an application configurator that allows some degree of control over the creation of non-salable data.

The application configurator is a software component that resides in the application domain and communicates configuration information about the application on its behalf such as the DSL specifications. Since this component operates without assistance from the application, it may exist in the form of an application library, or may be placed in the applications environment (via inheritance at execution time), or it can be implemented as a server process that proxies application information to the operating system as necessary.

A resource duplicator is a software component that fields requests for non-shareable resources and duplicates or virtualizes resources so that applications receive their own private copies and can co-exist transparently with multiple instances of the same application forged from the same virtual template. The resource duplicator also processes duplication rules fed by the application configurator or application snapshot/restore framework 200.

As used herein, non-salable data refers to any resource that is modified and globally visible to other application instances is non-salable (i.e. files). Process-related identifiers that are system-wide unique are also non-shareable since conflicts will arise if two instances use the same identifier at the same time (uniqueness is no longer preserved). References to unique resources by fixed handles (i.e. fixed TCP port numbers or IPC keys) are also not shareable. Memory pages that are specific to an application instance (i.e. the stack) are another example of a non-shareable resource. For illustrative purposes, examples of non-salable data include application config files that must be different per application instance as well as modified application data files if the application is not written to run multiple copies simultaneously. Other examples include stack memory segments or heap segments may also be non-salable data, shared memory keys that are a fixed value, usage of fixed well-known (to the application) TCP port numbers, and process identifiers (two distinct processes cannot share the same PID).

The snapshot virtual template is constructed automatically by dividing a snapshot process into shareable and non-shareable data. The knowledge of which system resources can be shared is encoded in the snapshot virtual templating framework itself. If an application has non-shareable internal resources (as opposed to system resources), it may not be possible to construct

a virtual-template for that application.

Snapshot virtual templates are node-independent as well as application-instance dependent. Snapshot virtual templates cannot be created for applications that use non-shareable physical devices. Snapshot virtual templates must save references to non-shareable resources in their pre- customized form, rather than their evaluated form. All access by an application to non-shareable resources must be via the operating system. Internal or implicit dependencies by the application itself cannot be virtually-templated. A snapshot virtual template may be created from an application instance that was originally forged from a different virtual template.

Snapshot virtual templating is an alternate method of creating an application instance.

10 The snapshot restore method described above requires creating unique instances of an application to create unique “snapshots” of that application. Virtual templating allows the creation of a generic application instance from which unique instances may be spawned. Every unique instance that is created from the original virtual template starts out as an exact copy (referred to herein as “clone”) but has been personalized just enough to make it a fully-functioning independent copy. Differences between copies may be due to the way resources are named or identified.

15 FIG. 7 illustrates in block diagram form the contents of a snapshot virtual template. The main components are resource name size, resource descriptor size, resource type, resource name, and resource data. Resource data includes many different types of information, as illustrated.

20 FIG. 8 describes the sequence of steps executed by the application snapshot/restore framework 200 to create a snapshot virtual template. As the application runs, every request for a new operating system resource (file, memory, semaphore, etc.) is checked for an existing rule. When the application is started under the virtual templating framework, a set of rules may be supplied at that time. The rule will state the type of resource, the type of access (i.e., create, read, 25 write, destroy, etc), and the action to be taken. If a rule is found (decision step 360), the rule is saved as part of the process state and recorded with the resource as auxiliary state at step 362. Rules may be added to the template that control the creation of application-instance specific resources. For example, environment variables or pathnames that incorporate an AID to differentiate and customize a particular resource among multiple instances. The following 30 syntax created for illustration purposes:

Define	<APPL-ID> as PROPERTY application-ID
REDIR PATH	"/user/app/config" to "/usr/app/<APPL-ID>/config"
SET ENV	"HOME" = "/usr/app/<APPL-ID>"

5 If rules are created, they should also be specified via the application configurator. If no rule is found, the resource is checked using a standard set of criteria that determine whether the resource needs to be abstracted or virtualized in order to be cloned at step 364. The criteria is again checked at steps 366, 370, 372, 378, 380 and 386. In most cases, no action is taken.

10 Resources are simply classified into their correct types so that when an instance is cloned the correct action can be taken. If the resource is shared, i.e. shared memory (decision step 366), the resource is marked as shared (step 368) so that during the subsequent snapshot all references to the shared object will be noted. If the resource can be modified (decision step 370), it must be isolated from the original during cloning so that the original remains untouched. If the resource is a large object and has a notion of an underlying object, such as i.e. mapped memory (decision step 372), it is marked for copy-on-write (step 374). Otherwise, the entire resource must be duplicated and marked accordingly (step 376). A resource is said to be systemwide unique if the identifier used to represent that resource cannot represent more than one instance of that resource at a single point in time on the same node or computer. If the resource is systemwide unique (decision step 378), and is exported as an external interface, as is the case when another client application that is not running on the platform has a dependency on the resource, such as a TCP port number (decision step 380), it isn't feasible to virtualize access to the resource, so it is marked to be multiplexed (step 382). Multiplexing allows multiple independent connections to transparently co-exist over the same transport using only a single identifiable resource. If it isn't externally exported, the resource is marked for virtualize at step 384. Continuing to decision step 25 386, if the resource is network-unique, it is marked for allocation at step 388. Control proceeds to step 390, where the resource request is processed. Steps 360 through 390 are repeated for every resource request occurring during application execution.

30 FIG. 9 illustrates the sequence of steps executed to perform cloning or replication of a process from a snapshot virtual template. This sequence of steps can be performed by a replication program that creates a new snapshot image from an existing template, or by the

restore driver 220. When an application instance is restored from a snapshot that is a virtual template, a new instance is automatically cloned from the template using the rules that were gathered during the creation of the template. For every resource included in the snapshot virtual template, rules for the resource and access type are looked up. Any resource that requires special handling as part of the templating effort has the rule described inside the snapshot template as part of the auxiliary state associated with the resource. If no rule is found (decision step 400), the resource is recreated using the existing saved information in the snapshot (step 402). Otherwise, if a resource is marked for duplicate (decision step 404), then a copy of the original resource is made at step 406. If a resource is marked for copy-on-write (decision step 408), then at step 410 a reference to the original underlying object (in the original template) is kept, and any modifications to the original force a copy-on-write action so that the modifications are kept in an application-instance private location and the two form a composite view that is visible to the application instance.

If a resource is marked for virtualization (decision step 412), the original resource is allocated or duplicated in blank form at step 414. At step 416, the resource is mapped dynamically to the new resource at run-time by binding the system resource to the saved resource in the snapshot image. If a resource is marked for multiplex (decision step 418), the original resource is duplicated and then spliced among other application instances that share it (step 420). If the resource is a network unique resource (decision step 422), a unique resource must be allocated (step 424) by communicating with another component of the network, i.e. network map or registry, that assigns a resource to this instance. Then this new resource is bound to the fixed resource that was saved in the virtual template (step 426), in a manner similar to virtualization.

### C. Virtual Resource ID Mapping

The present invention provides virtual mapping of system resource identifiers in use by a software application for the purpose of making the running state of an application node independent. By adding a layer of indirection between the application and the resource, new system resources are reallocated and then can be mapped to the application's existing resource requirements while it is running, without the application detecting a failure or change in resource handles.

This layer of indirection makes the application's system RID transparent to the application. RID's are usually numeric in form, but can also be alphanumeric. RID's are unique to a machine, and can be reused once all claims to a specific RID have been given up. Some examples of RID's include process ID's, shared memory ID's, and semaphore ID's. Only the virtual RID is visible to the application. Conversely, the virtual RID is transparent to the OS, and only the system RID is visible to the OS. Every application has a unique identifier that distinguishes it from every other running application. There exists a one to one mapping between the AID: resource type: virtual RID combination and the node ID: system RID. Virtual RID's are only required to be unique within their respective applications, along with their corresponding system RID's may be shared among multiple programs and processes that have the same AID. System RID's that have been virtualized are accessed through their virtual ID's to ensure consistent states.

AID's are farm-wide unique resources and are allocated atomically by the AID generator. Because in the present invention applications aren't uniquely bound to specific names, process ID's, machine hostnames or points in time, the AID is the sole, definitive reference to a running application and its processes. Typically, an AID is defined in reference to a logical task that the application is performing or by the logical user that is running the application.

Virtual resource mapping comprises several basic steps: application registration, allocation of the RID, and resolution of the RID. During registration of the application, the AID is derived if preallocated or the application existed previously, or it may be allocated dynamically by an AID generator. The AID is then made known for later use. Allocation of a RID happens when an application requests access to a system resource (new or existing) and the OS returns a handle to a resource in the form of a RID. The virtual resource layer intercepts the system returned RID, allocates a virtual counterpart by calling the resource specific resource allocator, establishes mapping between the two, and returns a new virtual RID to the application.

Resolution of a RID may occur in two different directions. A RID may be passed from the application to the OS, in which case the RID is mapped from virtual ID to system ID. Conversely, the RID may be passed from the OS to the application, in which case the transition is from system ID to virtual ID. Requests for translation are passed from the framework to the virtual RID translation unit and the corresponding mapping is returned once it has been fetched

from the appropriate translation table. Multiple translation tables may exist if there are multiple resource types.

FIG. 10 illustrates the steps executed to register an application. The AppShot harness 500 exists to aid in the creation of the appropriate runtime environment for the application prior to launching the application. The appshot harness 500 initializes the runtime environment for an application by first priming its own environment with the appropriate settings and then launching the application which inherits all these settings. The appshot harness 500 is provided because the application cannot be recompiled or rewritten to initialize its own environment. Some of the settings that are established by the appshot harness 500 include the AID, assigned process ID range, DSL specifications, application virtual ID's, and snapshot templating rules. The DSL specifications are registered as part of the environment. A process is an in-memory instantiation of a software program. Process<sub>x</sub> is the in-memory image of the AppShot harness 500 and process<sub>y</sub> is the in-memory image of the application. At step 510, the appshot harness 500 registers an AID a<sub>i</sub>, such as "dbserver." with the application snapshot/restore framework 200 within the OS kernel 206. The application snapshot/restore framework 200 then creates virtual translation tables 502 for the AID at step 512. Virtual translation tables 502 are data units that contain translation information related to RID's, such as AID's or process ID's, virtual RID's, and system RID's. Separate tables can be implemented per resource type or a table can be shared if a unique resource type is stored per table entry. A translation unit maps the system RID's to the virtual RID's by storing and fetching translation information in the appropriate translation table. Once the virtual translation tables 502 are created, Process<sub>x</sub> is linked to the AID a<sub>i</sub> at step 514. At this point, process<sub>y</sub> is created when the appshot harness 500 launches the application at step 516. Process<sub>y</sub> then inherits process<sub>x</sub>'s link to a<sub>i</sub> and its tables at step 518.

FIG. 11 shows the allocation of a virtual resource such as a semaphore in accordance with the present invention. At step 520, an application requests that a semaphore resource is allocated for its process<sub>y</sub>. In response (step 522), the application snapshot/restore framework 200 looks up the AID in memory 154, and returns AID a<sub>i</sub>. At step 524, in response to a request from the application snapshot/restore framework 200, the system semaphore pool returns semaphore s<sub>i</sub>. At step 526, the application snapshot/restore framework 200 scans the virtual resource translation table 502 for an available slot and allocates the virtual semaphore. At step 528, the

application snapshot/restore framework 200 inserts the translation  $s_i = a_i:v_3$  and the virtual resource translation table 502 now contains the mapping. At step 530, the virtual semaphore  $v_3$  is returned to the application.

FIG. 12 illustrates translation of a virtual resource to a system resource in accordance with the present invention. At step 532, the application calls the semaphore interface and supplies the virtual RID  $v_3$  to the application snapshot/restore framework 200. At step 534, the application snapshot/restore framework 200 looks up the AID for the calling application and returns  $a_i$ . At step 536, the application snapshot/restore framework 200 then looks up the translation for  $a_i:v_3$  in the virtual resource translation table 502, which returns  $s_i$  at step 538. At step 540, the OS semaphore implementation is achieved when the application snapshot/restore framework 200 forwards the application's request by substituting  $s_i$  for  $v_3$ .

FIG. 13 illustrates translation of a system resource to a virtual resource. Beginning at step 542, the application calls the semaphore interface and expects the RID as a result. At step 544, the application snapshot/restore framework 200 looks up the AID for the calling application and returns  $a_i$ . At step 546, the application snapshot/restore framework 200 forwards the application request to the OS semaphore implementation, which returns the system semaphore  $s_i$  at step 548. At step 550, the application snapshot/restore framework 200 then looks up the translation for  $a_i:s_i$  in the virtual resource translation table 502, which returns  $v_3$  at step 552. At step 554, the application snapshot/restore framework 200 returns the virtual RID  $v_3$  to the calling application.

FIG. 14 illustrates the logical sequence of steps executed to create the virtual translation table 502. Beginning at step 556, an attempt is made to register AID  $a_i$ . If the AID hasn't already been registered (decision step 558), a virtual resource translation table space for  $a_i$  is created in table 502 (step 560). Translation tables are then added for each type of resource associated with the application at step 562. At step 564, the process<sub>x</sub> is linked to  $a_i$ . At step 566, the process<sub>y</sub> is created and the application is launched. Steps 568 and 570 show the parallel paths of execution. The flow of control continues on to step 568 and halts shortly thereafter. The new flow of execution continues on from 566 to 570, where the process inherits context from the process at step 568.

FIG. 15 illustrates in greater detail the sequence of steps executed to translate a virtual

resource. For illustrative purposes, the resource in this example is a semaphore. Beginning at decision step 580, if an AID is found upon lookup, and the interface uses the RID as a parameter (decision step 582), the application snapshot/restore framework 200 performs a lookup of the translation for  $a_1:v_1$  at step 584. If a system resource  $s_1$  is found, the system RID is substituted for the virtual RID at step 586 and passed to the semaphore interface of the OS 206 (step 588). If the semaphore was not allocated by the semaphore interface (decision step 590), and the interface returns a semaphore (decision step 592) control proceeds to step 594 where a reverse lookup for the translation of the AID with a system RID is performed. The returned virtual RID is then substituted for the virtual ID at step 596. Returning to decision step 590, if a semaphore was allocated by the semaphore interface, control proceeds to step 598 where the virtual semaphore is allocated and a translation for  $v_2 = a_1:s_2$  is inserted into the translation table 502 at step 600.  $V_2$  is then substituted for  $s_2$  at step 602.

In another aspect, the present invention provides communication between at least two applications through virtual port multiplexing. The communication is achieved by accepting a connection from a second application on a first port and allocating a second port to receive the communication from the second application. Once the second port has been allocated the second port translation is recorded. The communication is sent to the first port from the second application and received on the second port. The communication is then delivered to a first application from the second port. In one embodiment the first application requests the communication from the first port and the first port is translated to determine the second port such that the communication is delivered to the first application in the step of delivering the communication to the first application.

In one embodiment, the communication is received on the first port following the step of sending the data to the first port, the first port is translated to determine the second port prior to the step of receiving the communication on the second port, and the step of receiving the communication on the second port includes queuing the communication on the second port from the first port.

In one embodiment, the second application requests to connect with the first port prior to the step of accepting the connection. Once the second port is allocated, the second port is negotiated including negotiating the second port between a first and second virtual port

multiplexer. Further, the second application is connected with the second port following the step of allocating the second port. The step of recording the translation including, first, recording the translation of the second port in association with the first application, and second, recording the translation of the second port in association with the second application.

5 The present invention also provides for a dynamic symbolic link (DSL) and the resolution of that DSL. The pathname of a first application is renamed to a target pathname, a variable within the target pathname, the first pathname is defined as a symbolic link and the symbolic link is associated with a virtual pathname. The method and apparatus further defines a specification is further defined that is associated with the virtual pathname including associating the variable 10 with the virtual pathname. In associating the symbolic link with the virtual pathname, a declaration is defined within the virtual pathname.

#### 60 D. VIRTUAL NETWORK IDENTITY

15 Virtualization of network identity is achieved by assigning a unique virtual IP address and virtual hostname to a group of processes that make up the application instance which the instance keeps throughout its execution. This virtual network identity stays with the application instance regardless of which node the application is running on. The framework, in essence, provides a mechanism to create this virtual network identity (VNI) around the application using 20 the virtual network parameters assigned to it. In one embodiment, the virtual network parameters include an IP address and hostname. The framework 200 ensures that the application's instance uses the virtual network parameters, transparently, so that it can be moved across machines, 25 without modifications to the application.

The virtual hostname resolves to the virtual IP address for both the applications registered with the VNI framework as well as those that are not registered. This may require configuration 30 of a name service or OS host configuration files. For example, if an application instance used a virtual IP address of 10.10.0.1 and a virtual hostname of host1055, the standard hostname to IP address resolution mechanisms (e.g. DNS or the /etc/hosts file) would have to be preconfigured to resolve a query of host1055 to IP address 10.10.0.1.

Any application configuration of addresses and hostnames uses the virtual hostname and 35 virtual IP address assigned to the instance. Virtualization of network identities is transparent to

the application running within a VNI. From the perspective of the application, the application is running on a single node which has an assigned IP address that corresponds to one of its network interfaces. The application requires no modifications to run in within the VNI framework.

A virtual address and virtual hostname are assigned to the application instance before the application is run. This virtual address may be statically preassigned or it can be dynamically assigned by an address resource manager. Registration of the virtual address and virtual hostname are made in a similar manner as the registration of applications using virtual ID's as described above in FIG. 10, including the steps of the appshot harness 500 registering the virtual IP address and virtual hostname with the application snapshot/restore framework 200, which in turn installs a virtual interface for the virtual IP address and records the IP address and hostname for the processes associated with the application. In this embodiment, the appshot harness 500 is a software module that drives the sequence to initialize the application snapshot/restore framework 200, register the virtual network parameters, install the virtual interface and run the application. The virtual IP address is unique to the application while the application is running. Similarly, the virtual hostname can be preassigned or dynamically assigned by an external entity or created using an algorithm based on the IP address to ensure uniqueness. When the application is to be run, the virtual IP address is allocated/installed as a virtual interface on the node. The virtual interface remains on the node as long as the application is running on that node. In one embodiment, a virtual network interface is a logical interface that allows a node to associate one or more IP addresses with existing physical or loopback network interfaces on the computer. This functionality is provided by some standard operating systems and allows the host to use one or more IP addresses as the local address for a single network interface.

FIG. 16 illustrates the resolution of the local hostname for an application running with the VNI. When the application is running, any calls the application makes to get the system hostname (hostname, gethostname, sysinfo, uname, etc..) at step 604 are intercepted by the application snapshot/restore framework 200. The application snapshot/restore framework 200 checks the process state associated with application for the virtual hostname (step 606), which returns the virtual hostname (step 608). At step 610, the application snapshot/restore framework 200 returns the virtual hostname. The application snapshot/restore framework 200 binds the application to its VNI by ensuring that the virtual address is chosen as the local

address, rather than a real address associated with one of the nodes' interfaces, whenever the application performs any network communications. Thus, when the application instance accepts a connection, or receives data, from a remote application, the local IP address chosen must be the virtual address. When connecting or sending data to a remote application, the local address again must be the virtual address for the application instance.

5 FIG. 17 illustrates how a server application running within a VNI is connected to by a client application. At step 612, the server application requests from the application snapshot/restore framework 200 to listen on port 9000, with a wildcard address. At step 614, the application snapshot/restore framework 200 requests the process virtual address, which is returned as 10.10.0.1 at step 616. Using the returned virtual address 10.10.0.1 as a local address, the application snapshot/restore framework 200 requests of the TCP/UDP module a listen/receive from port 9000 at step 618. The TCP/UDP module is the TCP/UDP transport layer module included in the operating system kernel. The client application then resolves the virtual hostname to 10.10.0.1 at step 620. The client application then connects to 10.10.0.1 port 9000 at step 622, completing the connection between client and server.

10 FIG. 18 illustrates how a client application running within a VNI connects with a server application. Beginning at step 624, the client application requests of application snapshot/restore framework 200 to connect/send to address 192.168.1.70 port 9001. The application snapshot/restore framework 200 requests the virtual address for the process associated with the client application at step 626, and is returned 10.10.0.1 at step 628. Using the virtual address returned at step 626 as the local address, the application snapshot/restore framework 200 requests to connect/send to 192.168.1.70:9001 of the TCP/UDP module at step 630. The TCP/UDP module performs a connect/send to at step 632, thus completing the connection between the client application having a VNI and the server application.

15 25 Having disclosed exemplary embodiments and the best mode, modifications and variations may be made to the disclosed embodiments while remaining within the scope of the present invention as defined by the following claims.